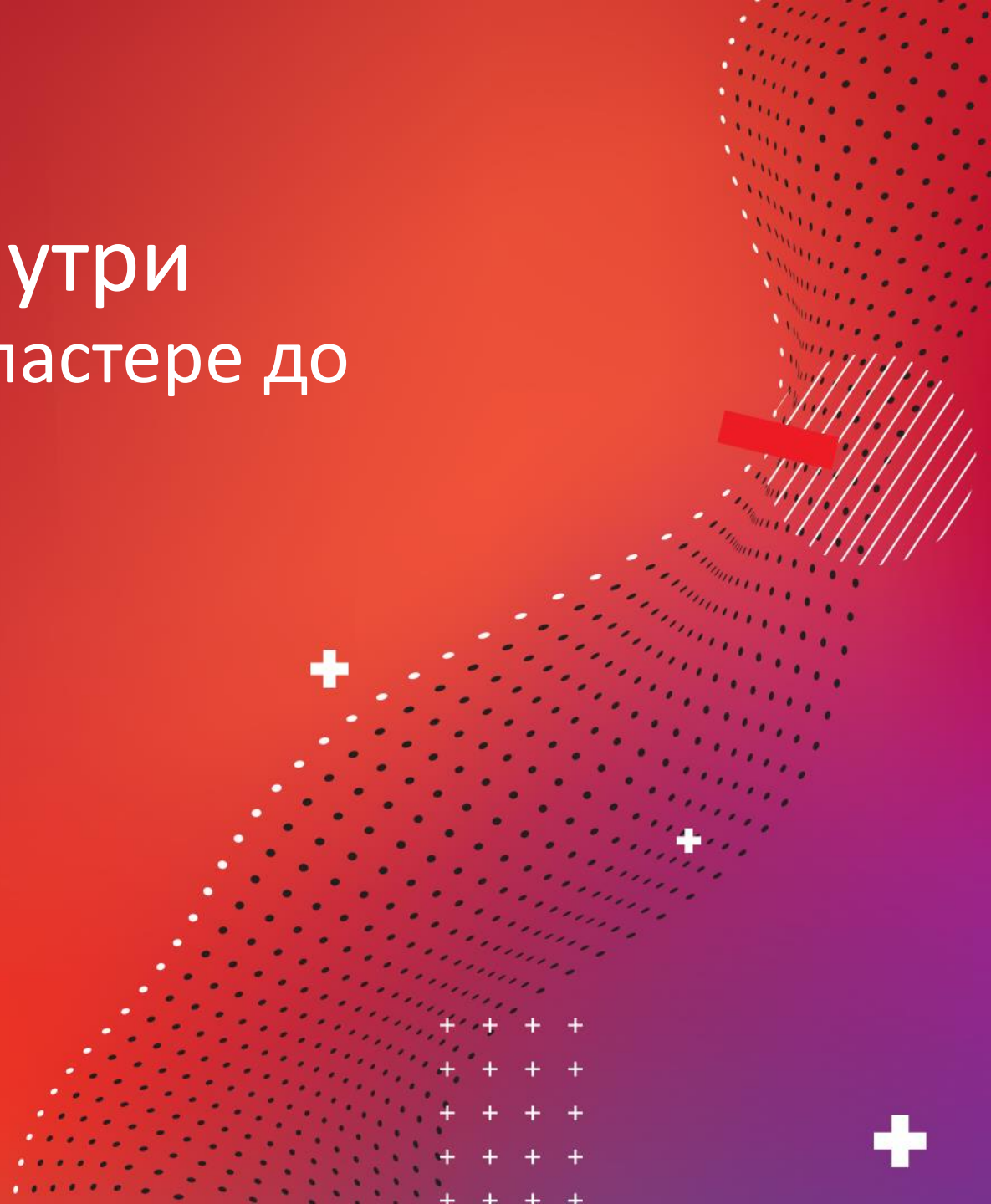# Apache Ignite, взгляд изнутри

от распределения данных в кластере до операции записи на диск

Дмитрий Павлов

HighLoad++
Весна 2021

# Dmitriy Pavlov

**VP, Apache Ignite**

Committer, PMC member

**Apache Training (incubating)**

Committer, PPMC member

**Customer Success Team lead**

Chief Technology Expert

Program committee member

14/18 years in Java development

dpavlov@apache.org

d-pavlov

# Disclaimer

This talk represents my own personal view and opinion.

It does not necessarily reflect the official stance of

The Apache Software Foundation/SberTech/

any company/any other entity the author might be affiliated with at the moment of presenting or in the past.

# Apache Ignite
# Intro

is a Distributed Database
for High-Performance Computing
with In-Memory Speed

1 Ignite cluster – N Caches

Cache - key-value storage
- put(k,v)
- v=get(k)

In Memory Data Grid – yes
In Memory Database - yes
SQL support - yes
SQL database – not fully

**IN-MEMORY DATA GRID**

= Stores data in-memory
 (data grid)

+ Compute - code goes to data
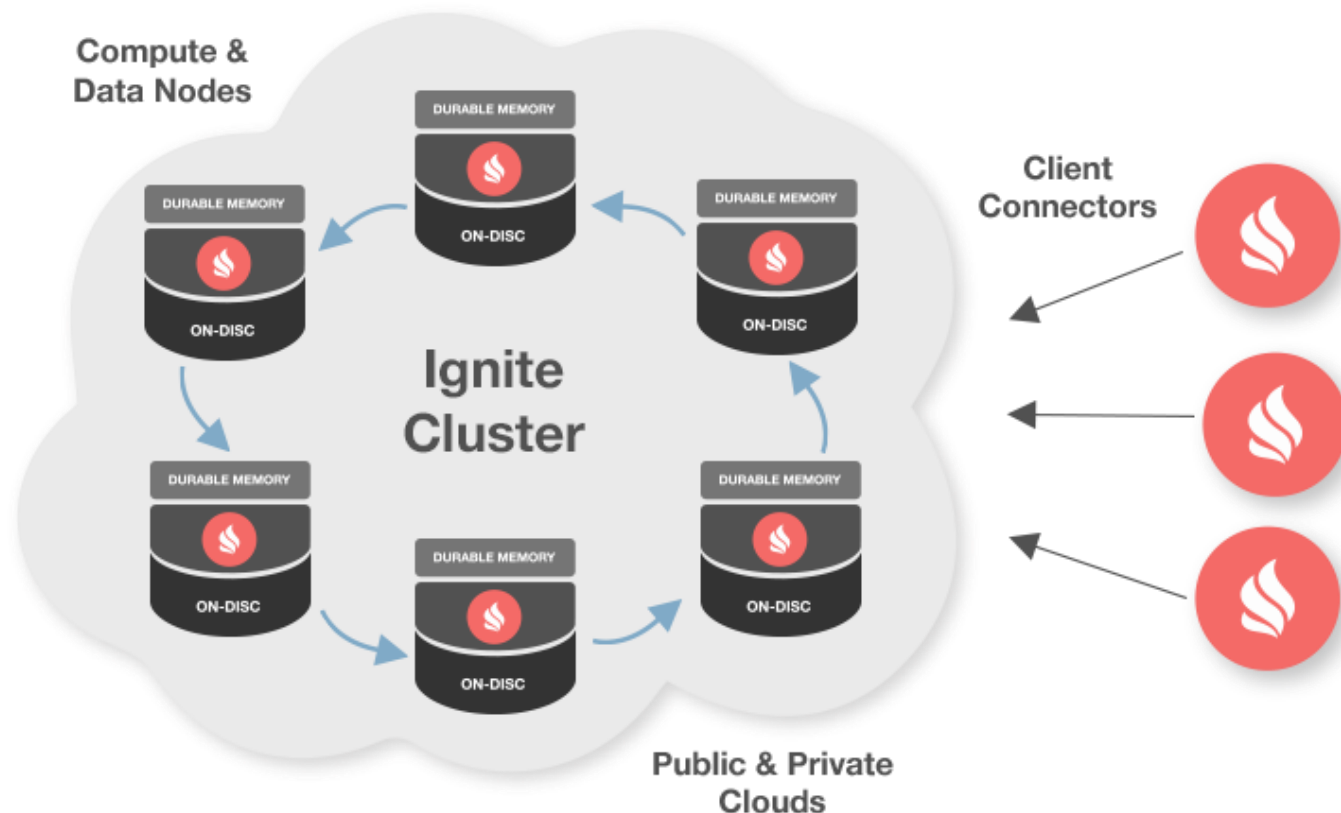(compute grid)

**DISTRIBUTED DATABASE**

• Memory-centric database, since V2.1

• Scalable: Each node stores only it's own data part

Ignite can be used in combined mode (part in-memory, part - persisted)

External data source (DB, REST, other)
• Yes - Cache
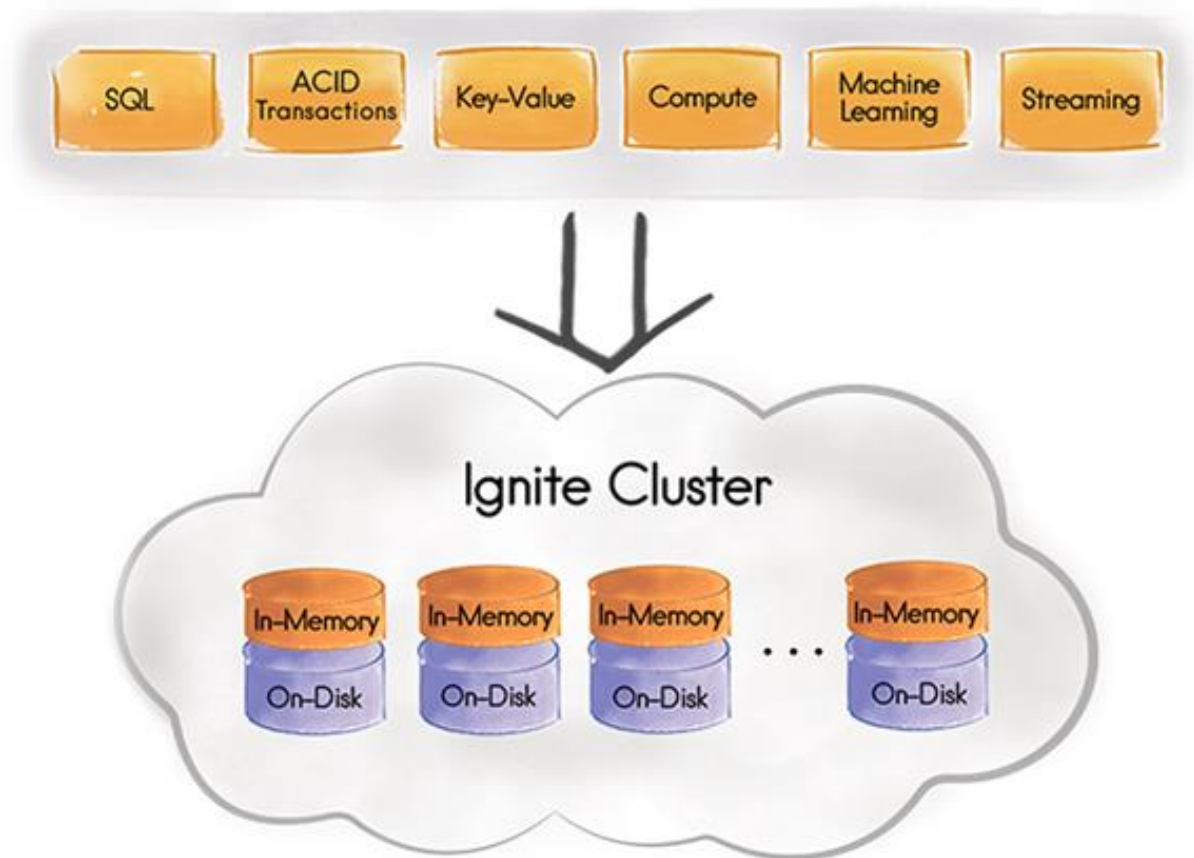• No - Ignite is Primary storage

# Ignite Cluster



- Servers – can store data
- Clients (thick) – put/get
- + Thin clients
- data distribution is handled by affinity function ([Rendezvous Hashing](#))

# Apache Ignite Native Persistence

is a distributed ACID and
SQL-compliant disk store

All data is on-disk,
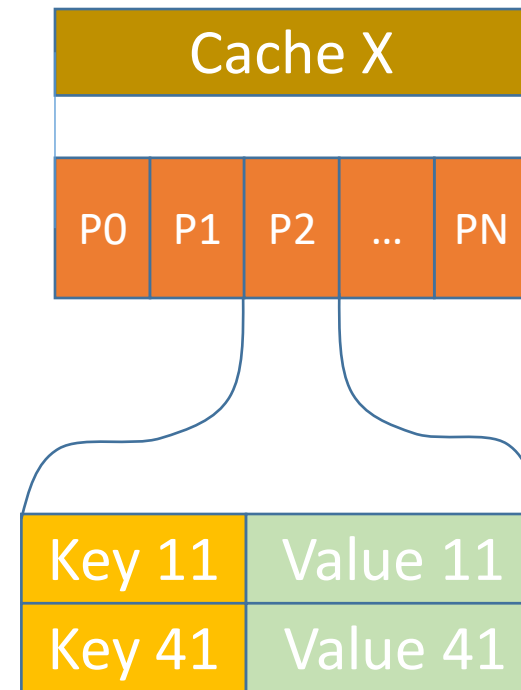part of the data is in-memory

Apache Ignite =
Speed/Scale



Each node has its own local storage
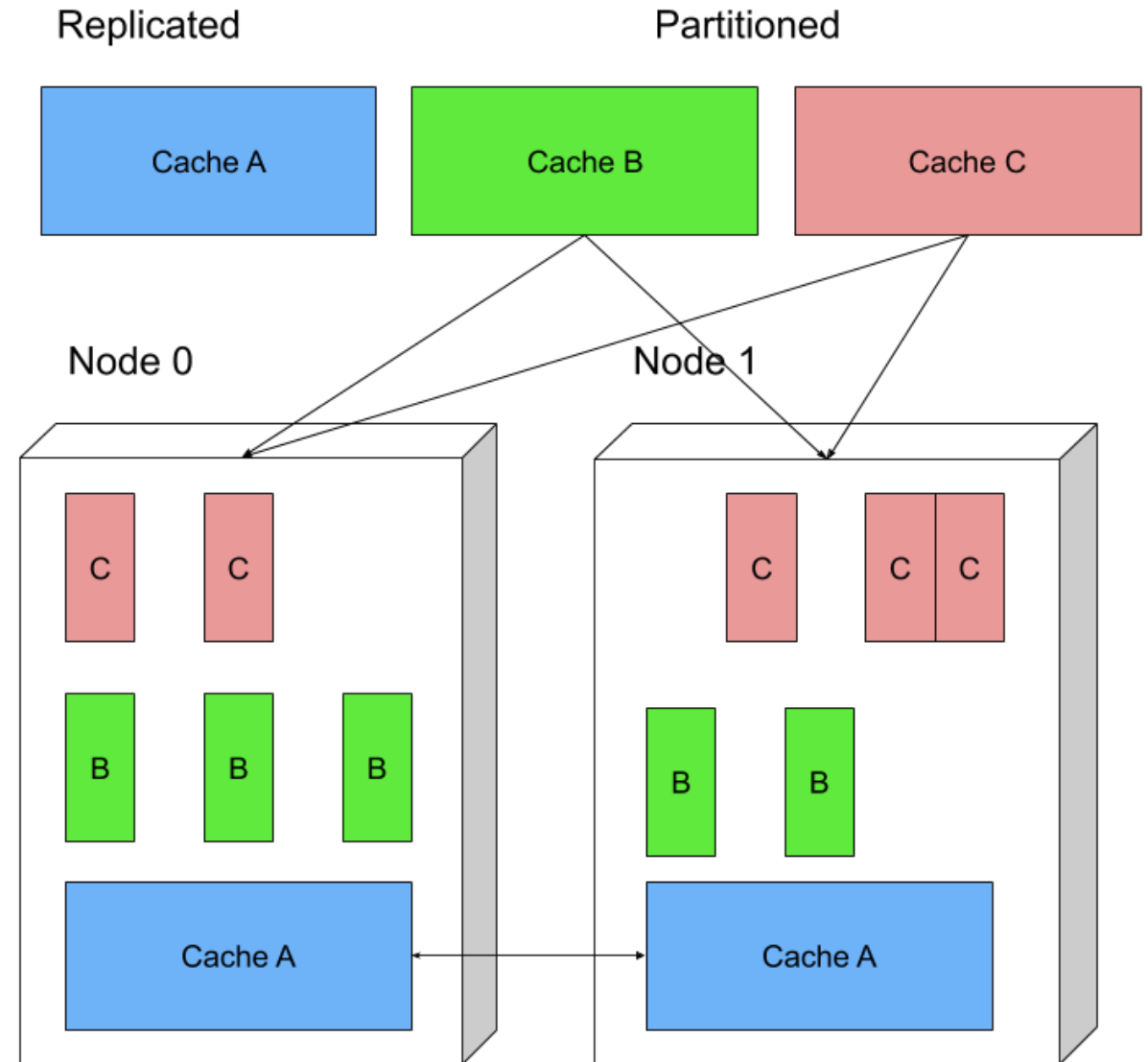
# Apache Ignite
## Cache

# Ignite Cache

- Cache main feature: K->V

- Cache ~ Java.util.Map

- JSR-107, JCache API

- Cache ~ Table

- Cache usually stores 1 business entry

- Entry = K + V

- 1 Cache – N partitions

- K -> hash(K) ->partition -> node

| Cache X | | | | |
|---|---|---|---|---|
| P0 | P1 | P2 | … | PN |

| Key 11 | Value 11 |
|---|---|
| Key 41 | Value 41 |

# Ignite Cache types

- Partitioned and Replicated
  - Replicated - Cache A
  - Partitioned - Cache B & C

- Replicated, use case:
rare write, often - read, e.g. dictionary

- Partitioned – most common
1024 default

- Backups 0,1…

# AGENDA

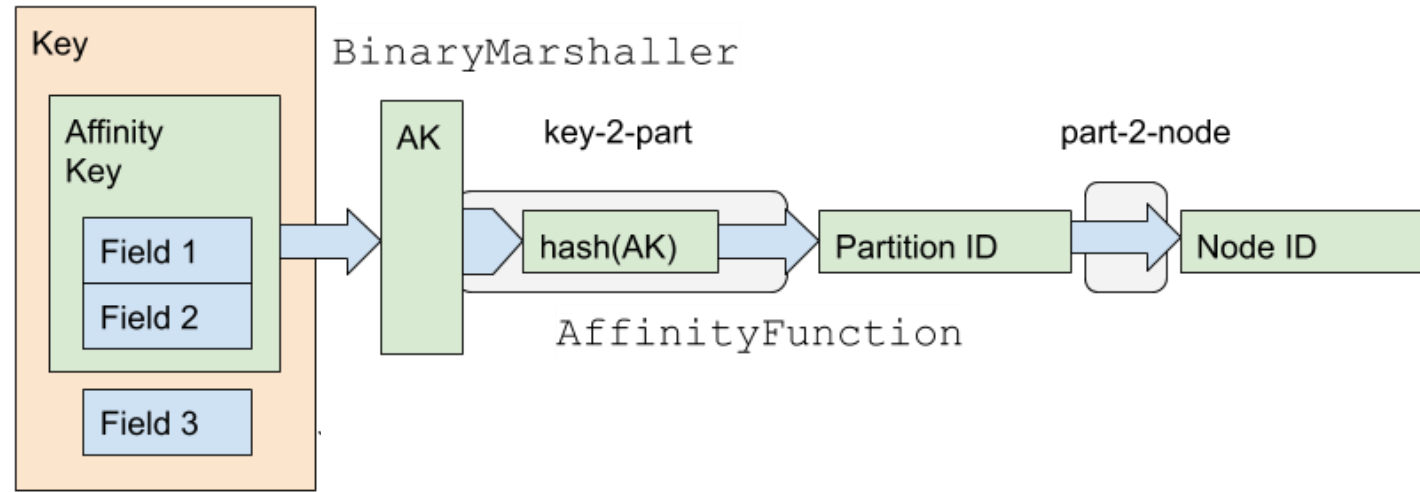What is Apache Ignite and caches

How to find an ideal match node

How to find some entry and update it

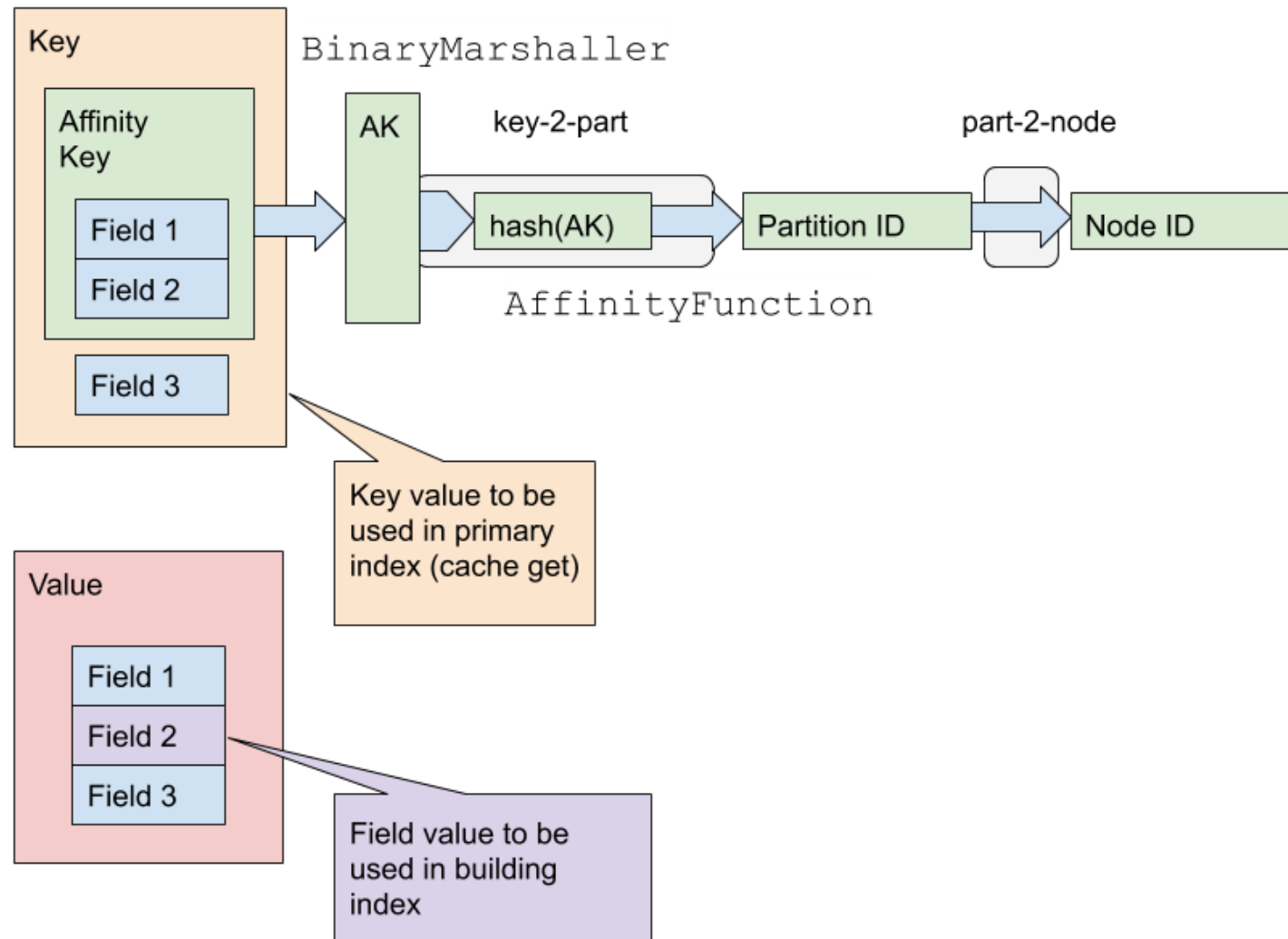How the storage and redo log are organized

Over speed protection

# Data location in the cluster
# Affinity

# Cache Key to node: Affinity Function



Affinity key
- is a key which will be used to determine a node

# Cache Key to node: Affinity Function



Key (full) is needed for primary index inside partition.

Key and indexed fields serialized in the similar manner

# Affinity function alternatives

- Naive : targetNode = K.hashCode() mod nodeCount

- Consistent Hashing (no more used in Ignite)

  https://en.wikipedia.org/wiki/Consistent_hashing

  http://theory.stanford.edu/~tim/s16/l/l1.pdf

- Rendezvous – used by Ignite

  https://en.wikipedia.org/wiki/Rendezvous_hashing

  http://www.eecs.umich.edu/techreports/cse/96/CSE-TR-316-96.pdf

# Rendezvous or Highest Random Weight (HRW)

Minimizes rebalancing on nodes set changes (leave node, join node)


- Naive approach node = K.hashCode() %  nodes

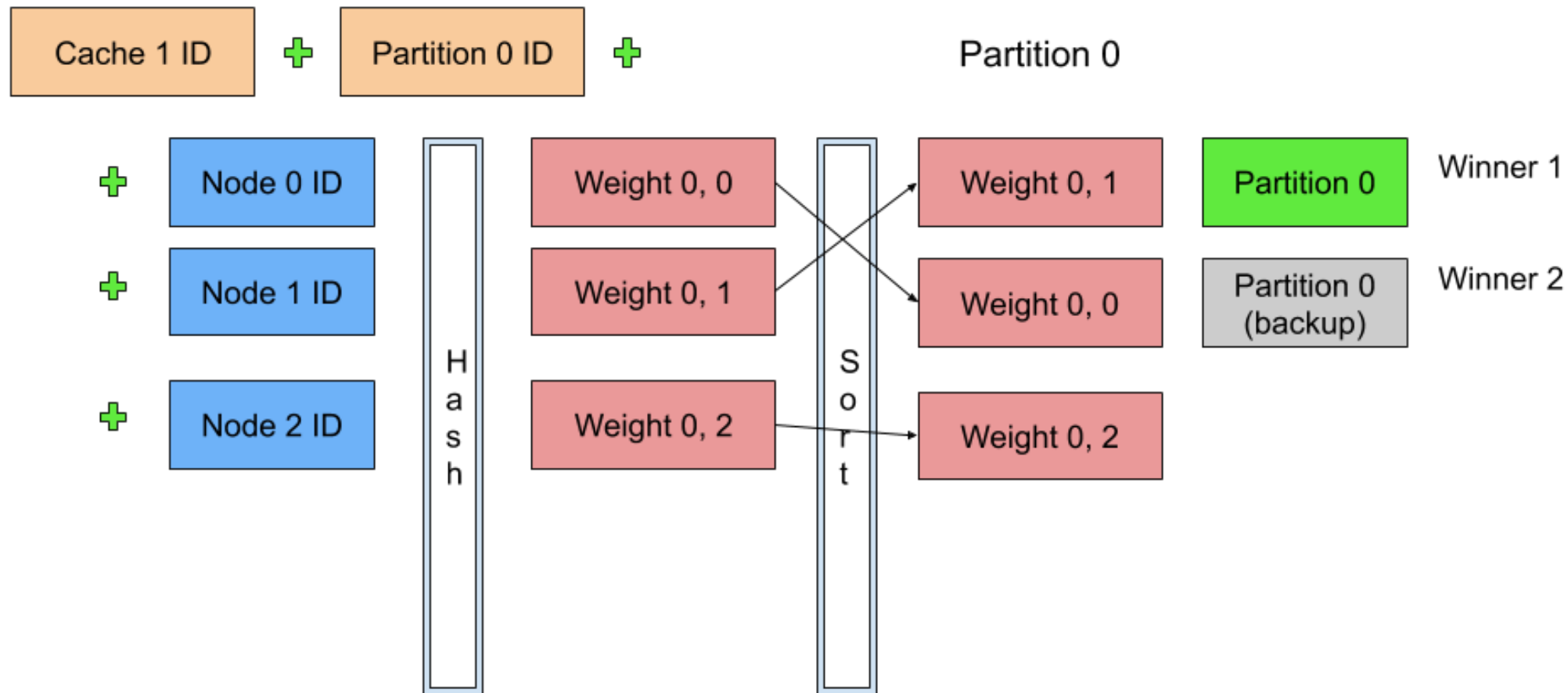- Imagine node add event

Consider case 1 cache 1024 partitions

```
012345… %3 = 012012012…
012345… %4 = 012301230…
```

- Almost all partitions should migrate

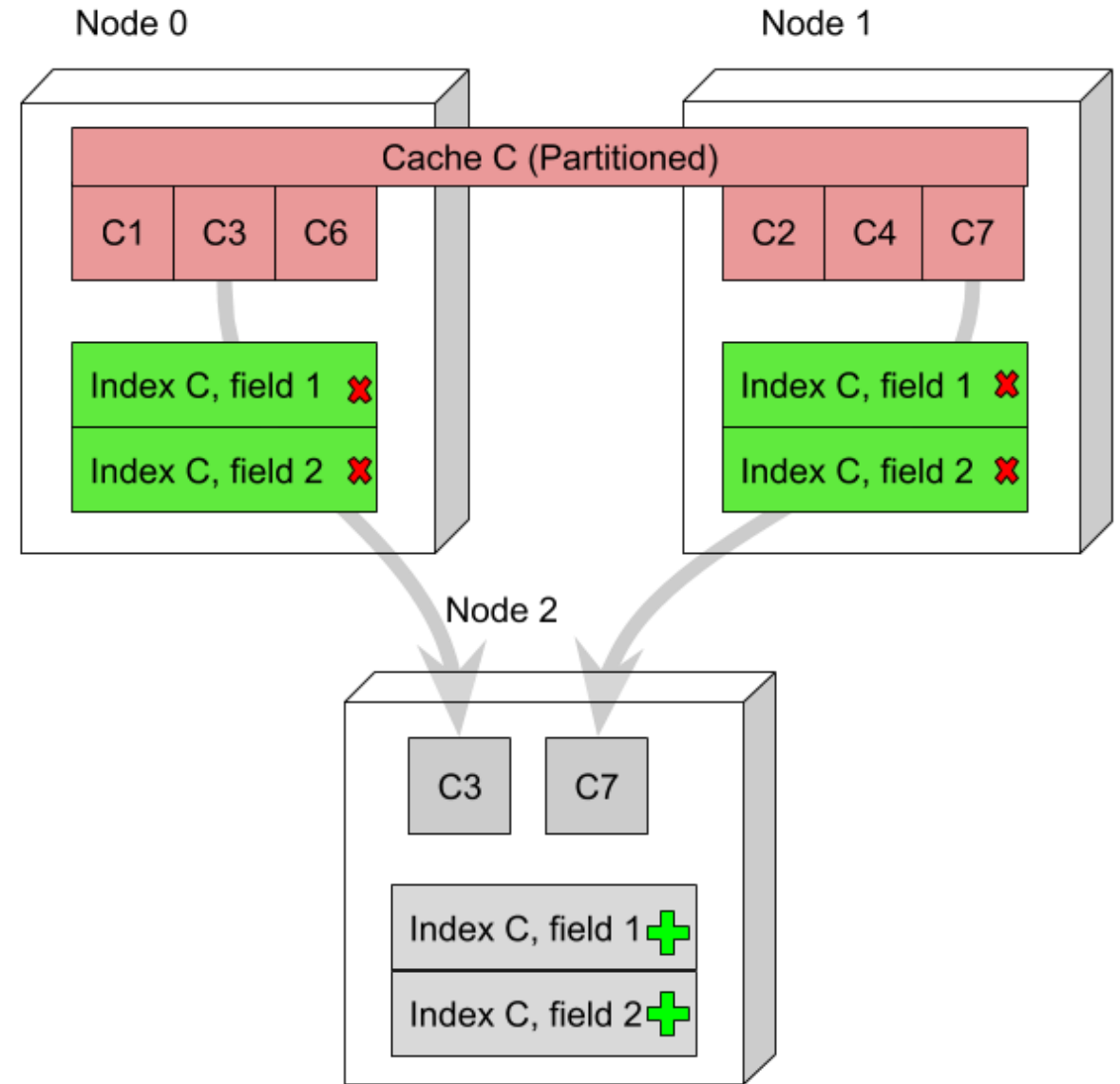| Cache 0 |
|---------|

| Partition 0 |
|-------------|

Affinity gives ideal topology.

(Cache partition target node.
Actual data may be on the way )

- Rebalancing: Moving data from one node to another

- Actual get should go to old node until rebalancing finished

- GridDhtPartitionFullMap (simply - node2part)

- Indexes affected during rebalancing

## AGENDA

What is Apache Ignite and caches

How to find an ideal match node

How to find some entry and update it
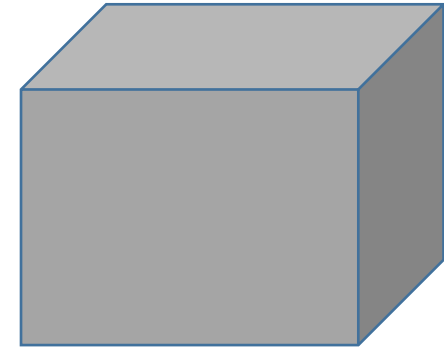
How the storage and redo log are organized

Over speed protection

LOCAL node:
Locating a key

# Let's move closer to the disk

**Node 0**

- One node

- One cache

- One partition

- One key

- We have some value to find

- Key is serialized by a marshaller (usually – binary)

- Split value (and key) - to chunks/blocks/pages

Cache 0

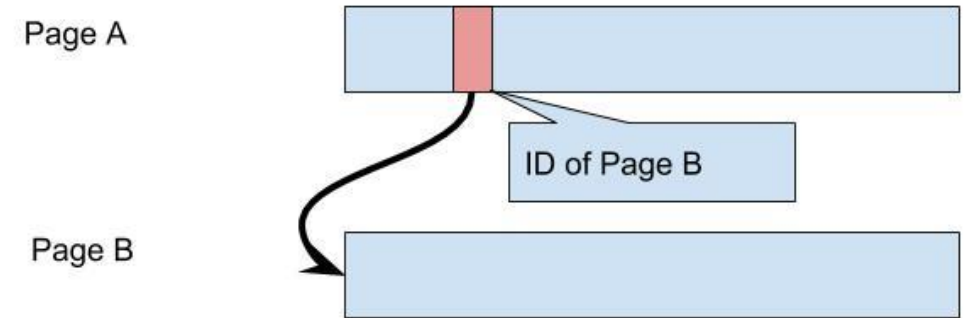Partition 0

Key 0  Value 0

HighLoad++
Весна 2021

All HDDs are block devices

Durable (Page) memory
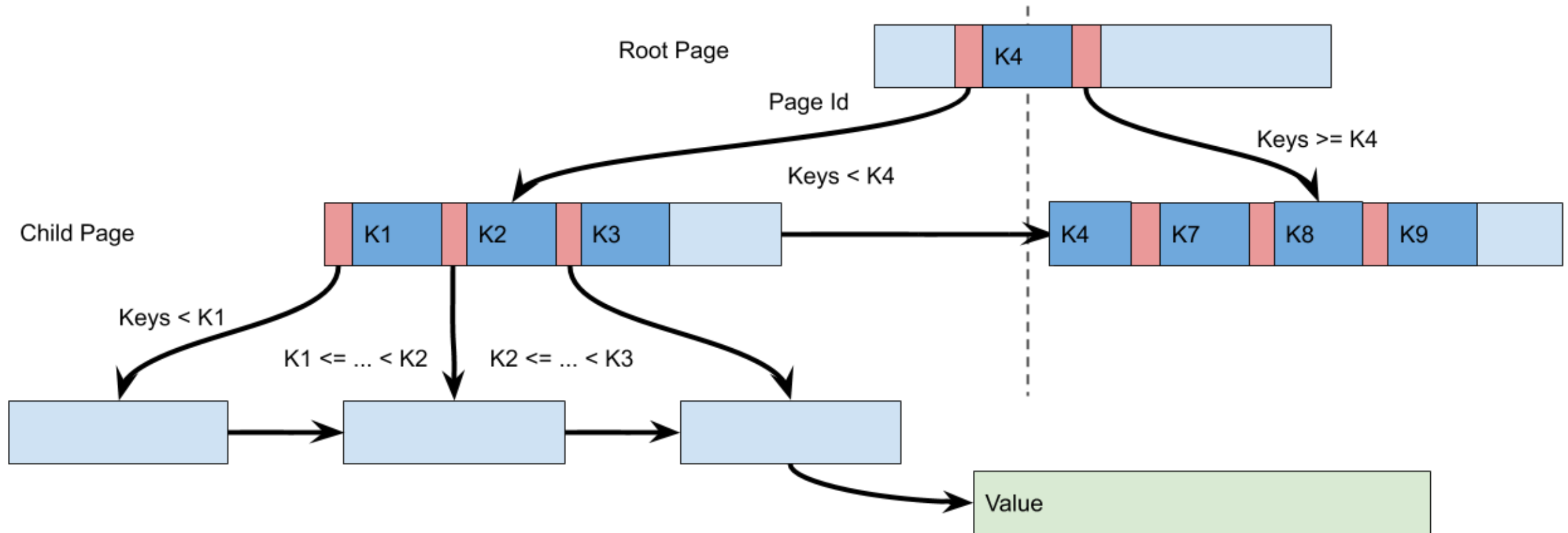
# Pages Identification

- 4k

- Index – int: 0,1,..

- (+) Partition ID = Page Id


- Links ~ "Pointers"

- Links survives memory-HDD-memory
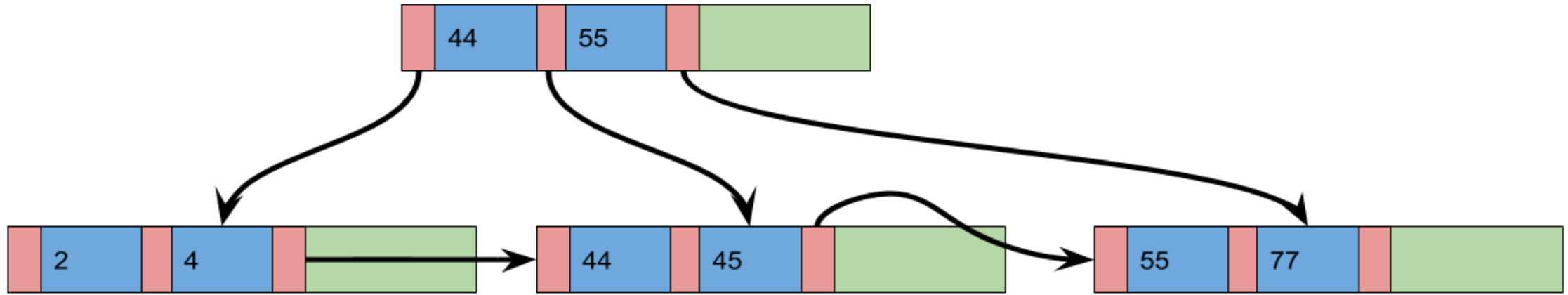
      (not depend on real address)

# How to find a key in our pages

- B+ Tree

- Read optimized

- ~ Linked list (from pages) with top levels

- Order of iteration - preferable for range lookups

- PK/Primary Index
  - for each partition
  - Key Hash based
  - Key value compared (collision resistant)
- Secondary index: value or value start in Index page.

# B+ Tree



Root Page

K4

Page Id

Keys < K4

Keys >= K4

Child Page

K1   K2   K3

K4   K7   K8   K9

Keys < K1

K1 <= ... < K2

K2 <= ... < K3

Value

# B+ Tree



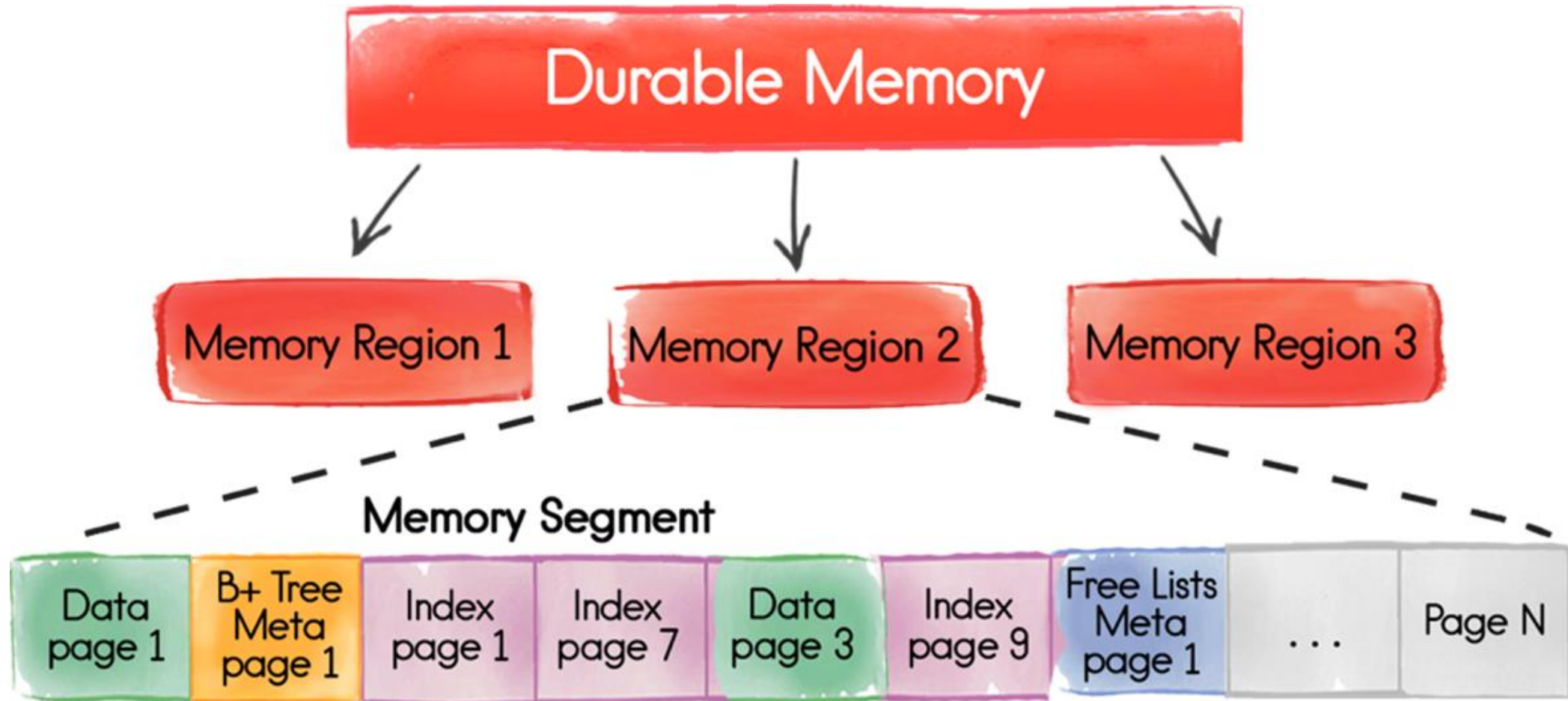https://www.cs.usfca.edu/~galles/visualization/BPlusTree.html

# LOCAL node:
# Store Values

# Pages are allocated within region randomly

# Data Page structure

| Page | | |
|---|---|---|
| | Page Header | Page Data |

| Data Page | | |
|---|---|---|
| | Page Header | Data Page Data |

**Entries**

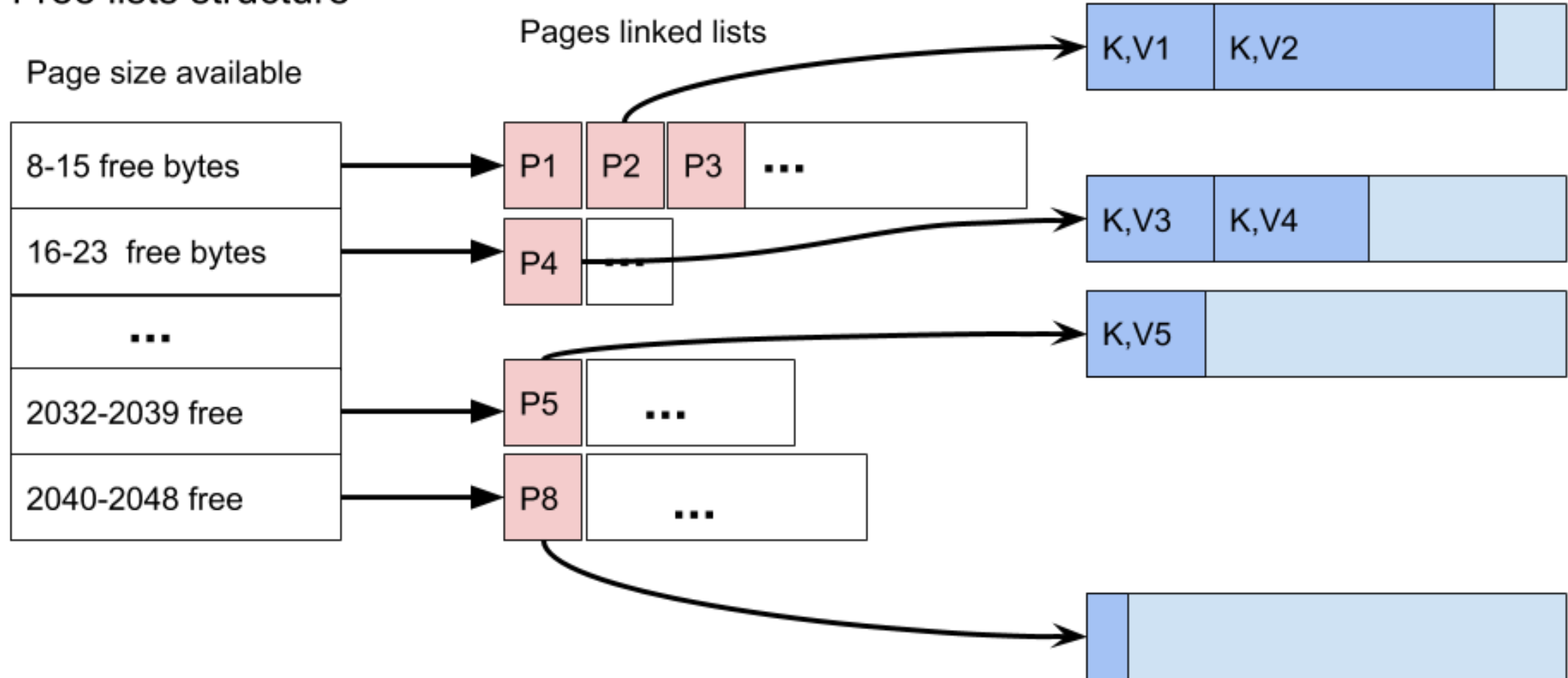| It1 | It2 | It3 | Free space | | K,V3 | K,V2 | K,V1 |
|---|---|---|---|---|---|---|---|

K&V locator (in local node): FullPageId + item

# Find suitable page for insertion data

Free lists structure

Page size available

Pages linked lists

| 8-15 free bytes | → P1 P2 P3 ... |
| 16-23 free bytes | → P4 ... |
| ... | |
| 2032-2039 free | → P5 ... |
| 2040-2048 free | → P8 ... |

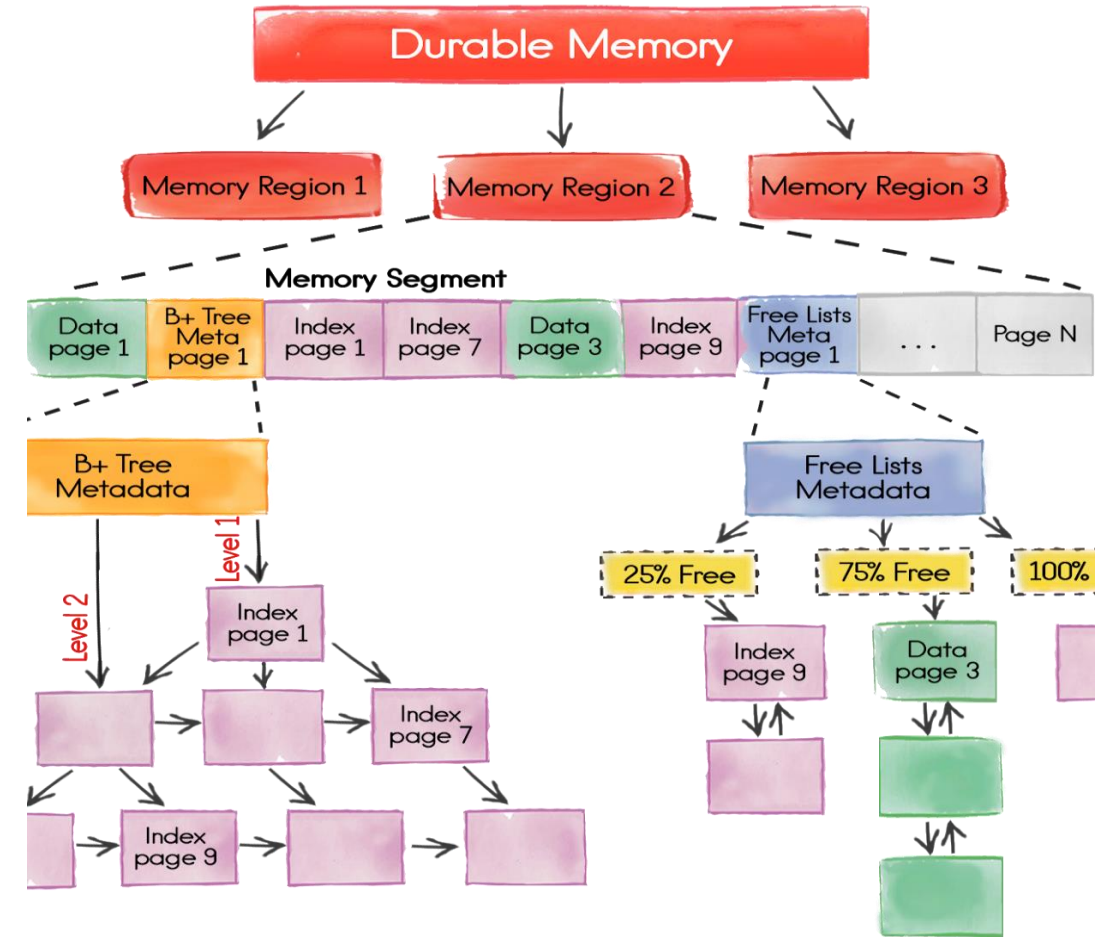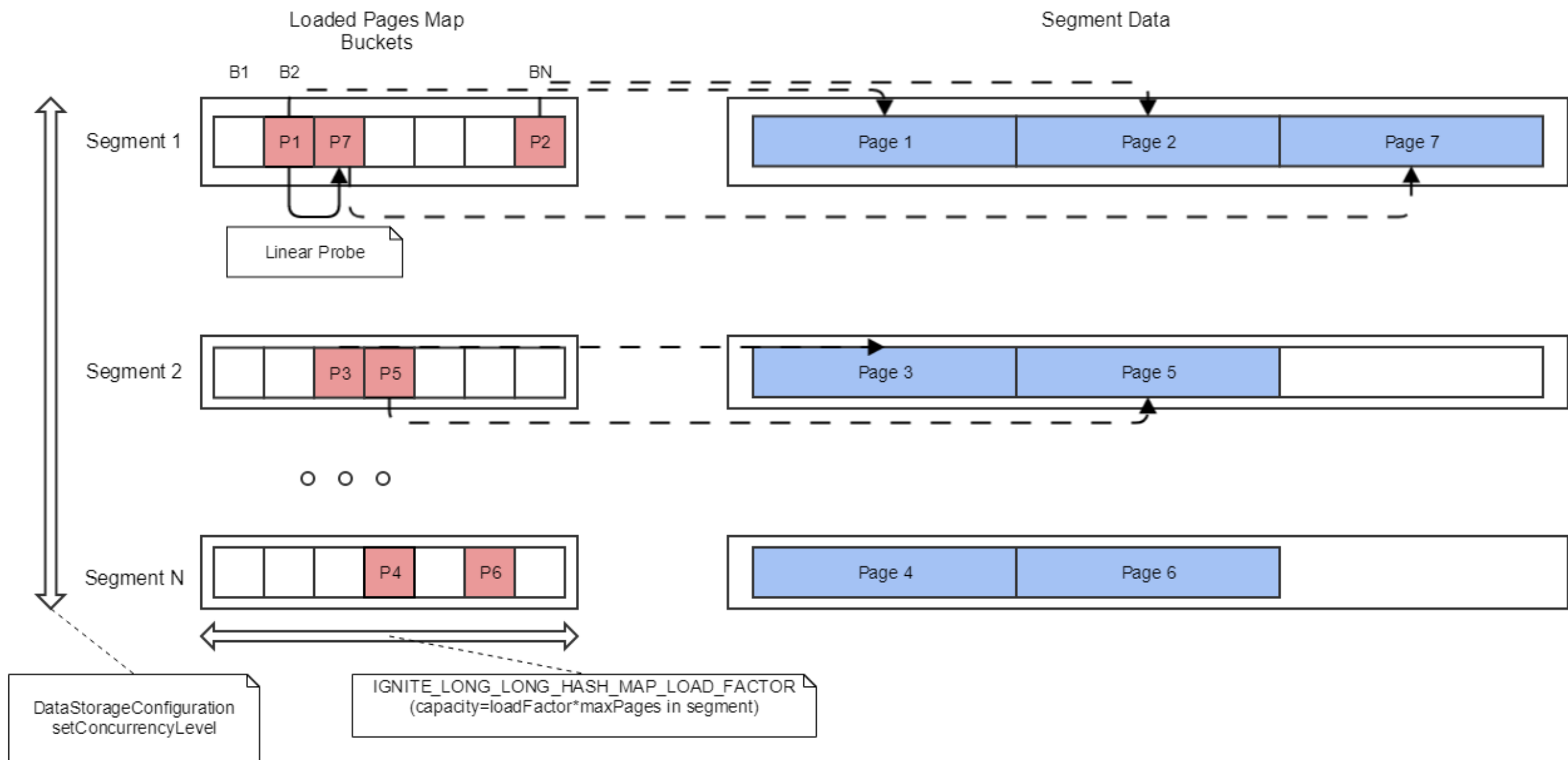K,V1 | K,V2

K,V3 | K,V4

K,V5

# Long objects

# RAM structure overview

- Regions – as configured
- Region has segments (depending on CPU count by default)
- Segment = set of RAM pages
- Pages has types and different formats
- Pages are linked between each other (Cross segment)

Loaded Pages Map Buckets

Segment Data

B1  B2  BN

Segment 1

| | P1 | P7 | | | | P2 |

Linear Probe

Page 1  Page 2  Page 7

Segment 2

| | | P3 | P5 | | | |

Page 3  Page 5

Segment N

| | | | P4 | | P6 | |

Page 4  Page 6

DataStorageConfiguration
setConcurrencyLevel

IGNITE_LONG_LONG_HASH_MAP_LOAD_FACTOR
(capacity=loadFactor*maxPages in segment)

HighLoad++
Весна 2021

LOCAL node:
A Page Modify and Write

# Page R/W operations in RAM

Page ID->real address in memory

-1 atomic operation: resolution of ID to address & lock page

Case of setting a field in the page:

o `AbstractDataPageIO#setFirstEntryOffset`

o `PageUtils#putShort`

o `GridUnsafe#putShort(long, short)`

o `sun.misc.Unsafe#putShort(long, short)`

# Read-Write (Memory & disc)

```
// Ignite classes: * FileIO

RandomAccessFileIO#write(ByteBuffer, long)

FileChannelImpl#write(ByteBuffer, long)

IOUtil#write(FileDescriptor, ByteBuffer, long,

                NativeDispatcher)

OUtil#writeFromNativeBuffer

NativeDispatcher#pwrite(…,

                (DirectBuffer)var1).address() + (long)var5
```

# IO Util implementation

```
if (var1 instanceof DirectBuffer) {
    return
        writeFromNativeBuffer(var0,var1, var2, var4);
}  else {
    …
    ByteBuffer var8
        = Util.getTemporaryDirectBuffer(var7)
}
```

## AGENDA

What is Apache Ignite and caches

How to find an ideal match node

How to find some entry and update it

✚ How the storage and redo log are organized ✚

Over speed protection

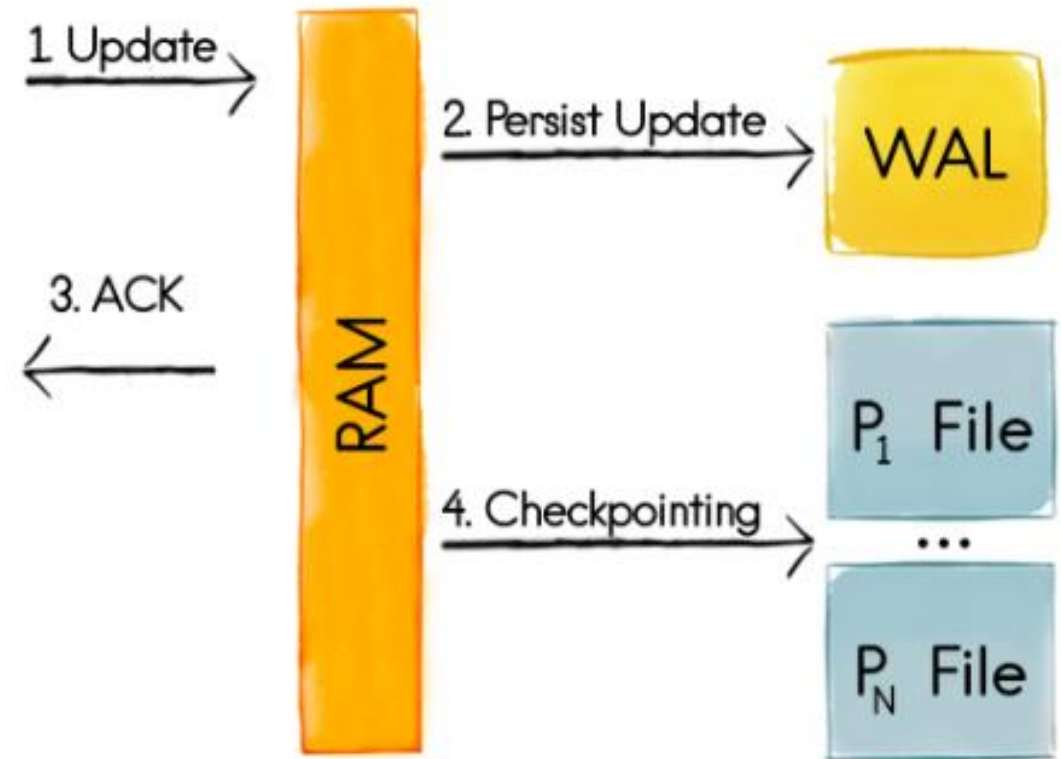WAL

# Write is not immediate after update

Data comes to
(1) In memory
Pages become dirty

(2) Write Ahead Log
(3) Update/TX completed

(4) Checkpointing = updating page store
files = Background process



1 Update

2. Persist Update → WAL

3. ACK

RAM

4. Checkpointing

$P_1$ File

...

$P_N$ File

# WAL

- WAL =  **A**CI**D** – A&D, properties

[https://en.wikipedia.org/wiki/Write-ahead_logging](https://en.wikipedia.org/wiki/Write-ahead_logging)

- Both logical
    - Set user.lastSeen=…
- And physical
    - Change page PageID=…, at offset 4 to NNNNNNNN
-

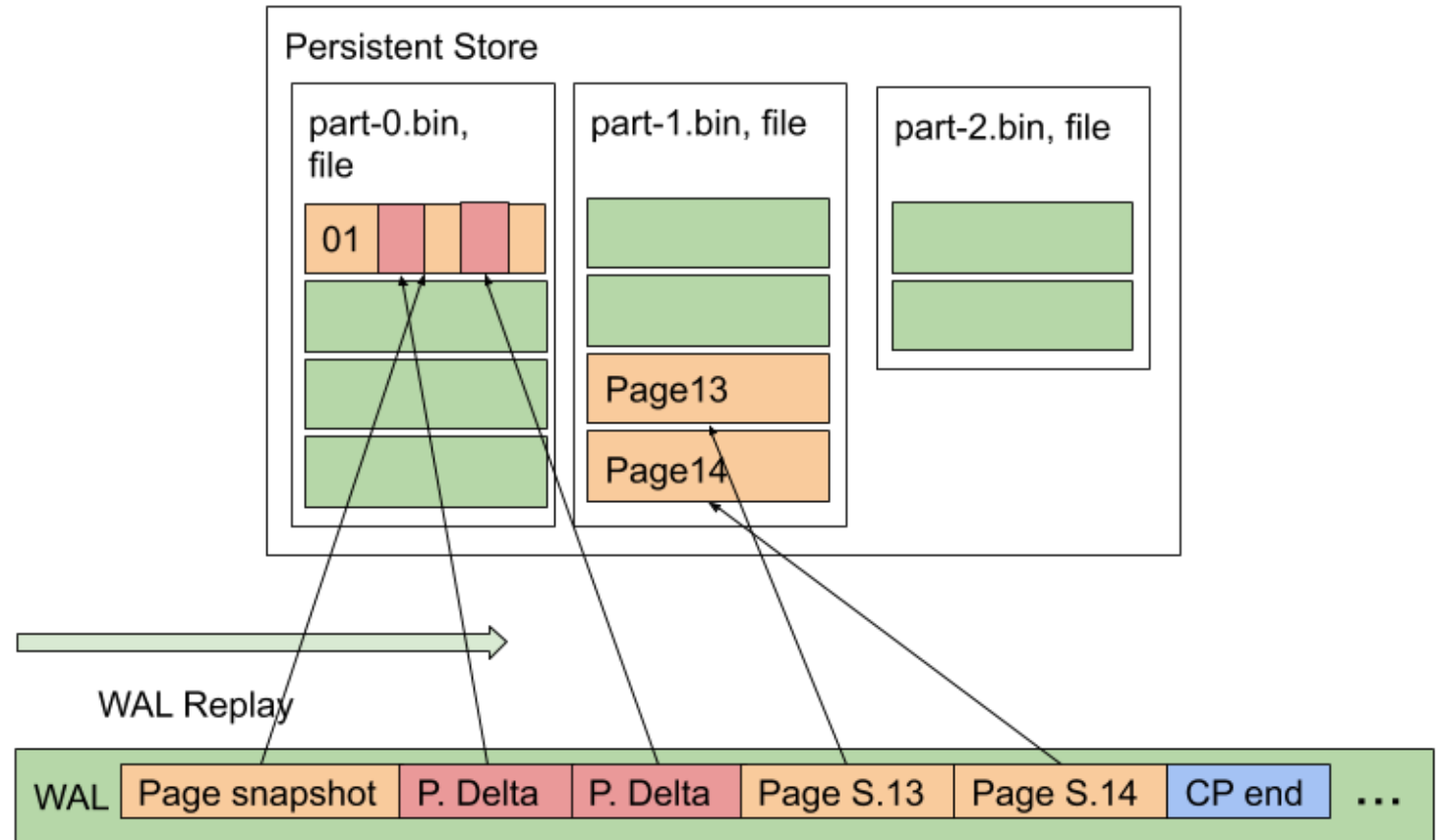# Operation example

# Why 2 types of records

- After checkpoint

- -1 CP history for data (logical)

- In the middle of checkpoint

- -2 CP for structure (physical)

- -1 CP for data (logical)
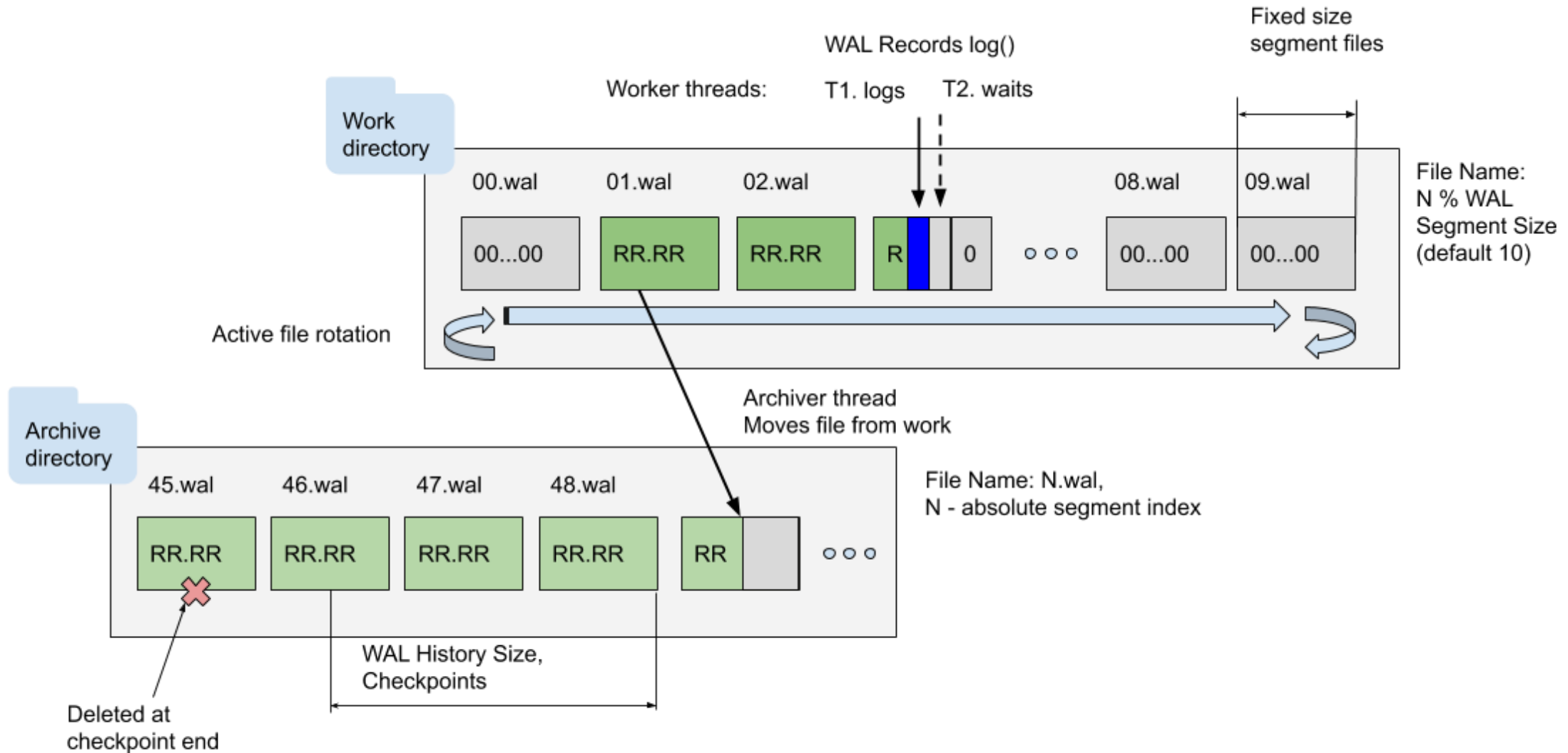
# More than 1 record for 1 put

- 1 logical, 1 physical

- Field changes its length + free lists updates

- Field update for 2+ pages for long objects

- Indexed field updated, need to add and remove index B+Tree nodes

- Updates in index may require Split-Merge operations for pages (2-3 page affected)

- Ignite tracks modifications, so Tracking page will be updated


- BUT: Latest Apache Ignite can share byte payload between records.

# Sync strategies

- FSYNC – any case survive – OS crash and power

- LOG_ONLY – give the data to the OS - process crash

- BACKGROUND – by timer, some records may be lost


- Actual WAL is not one file

- Set of files = segments

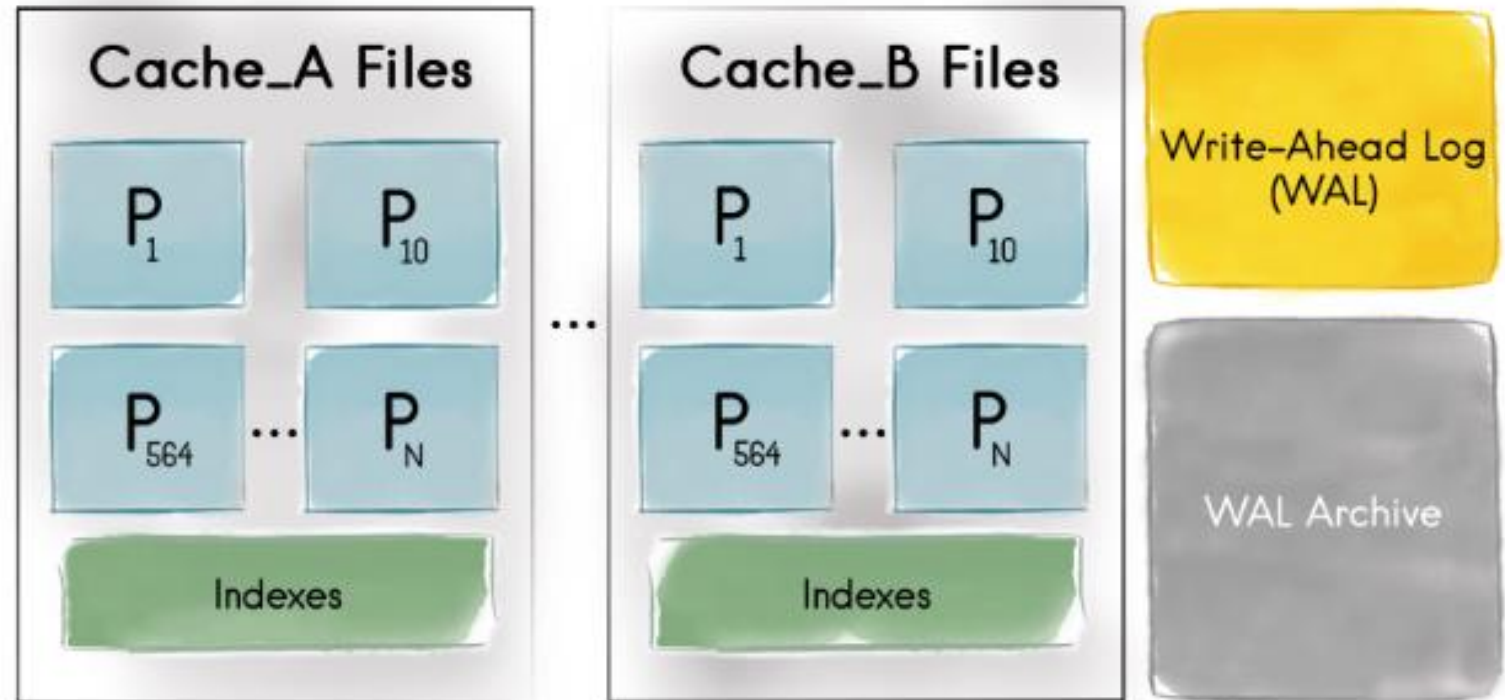- active, ready to be filled – Work
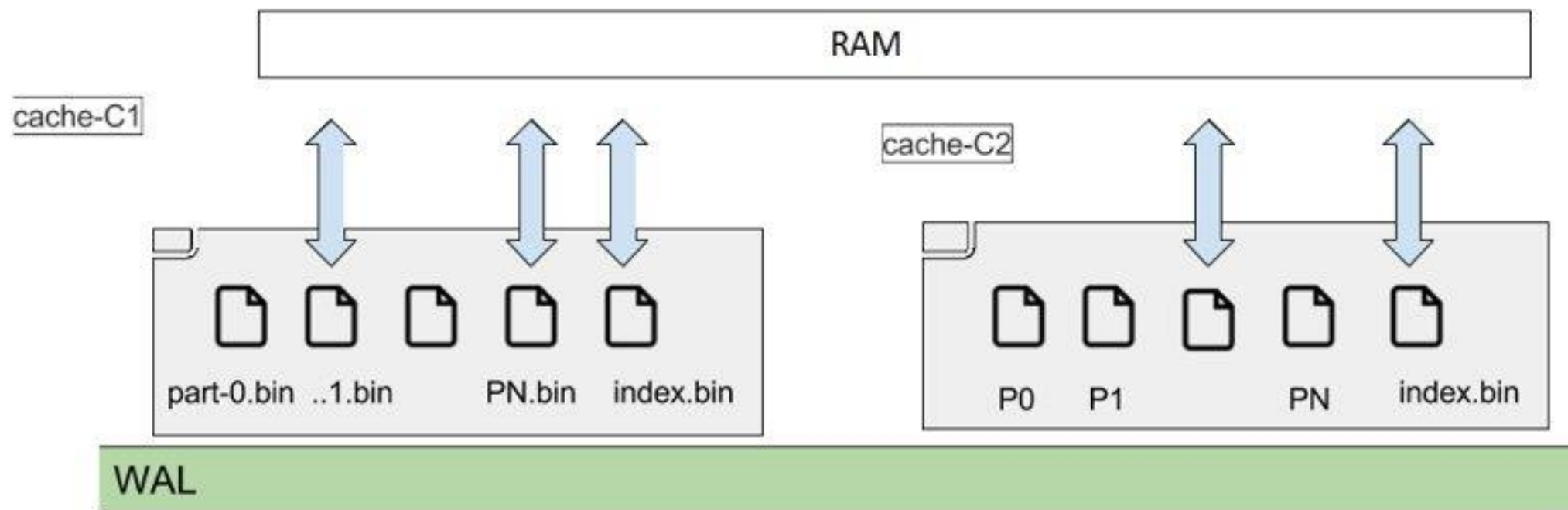
- Finalized – archive

# Actual WAL structure

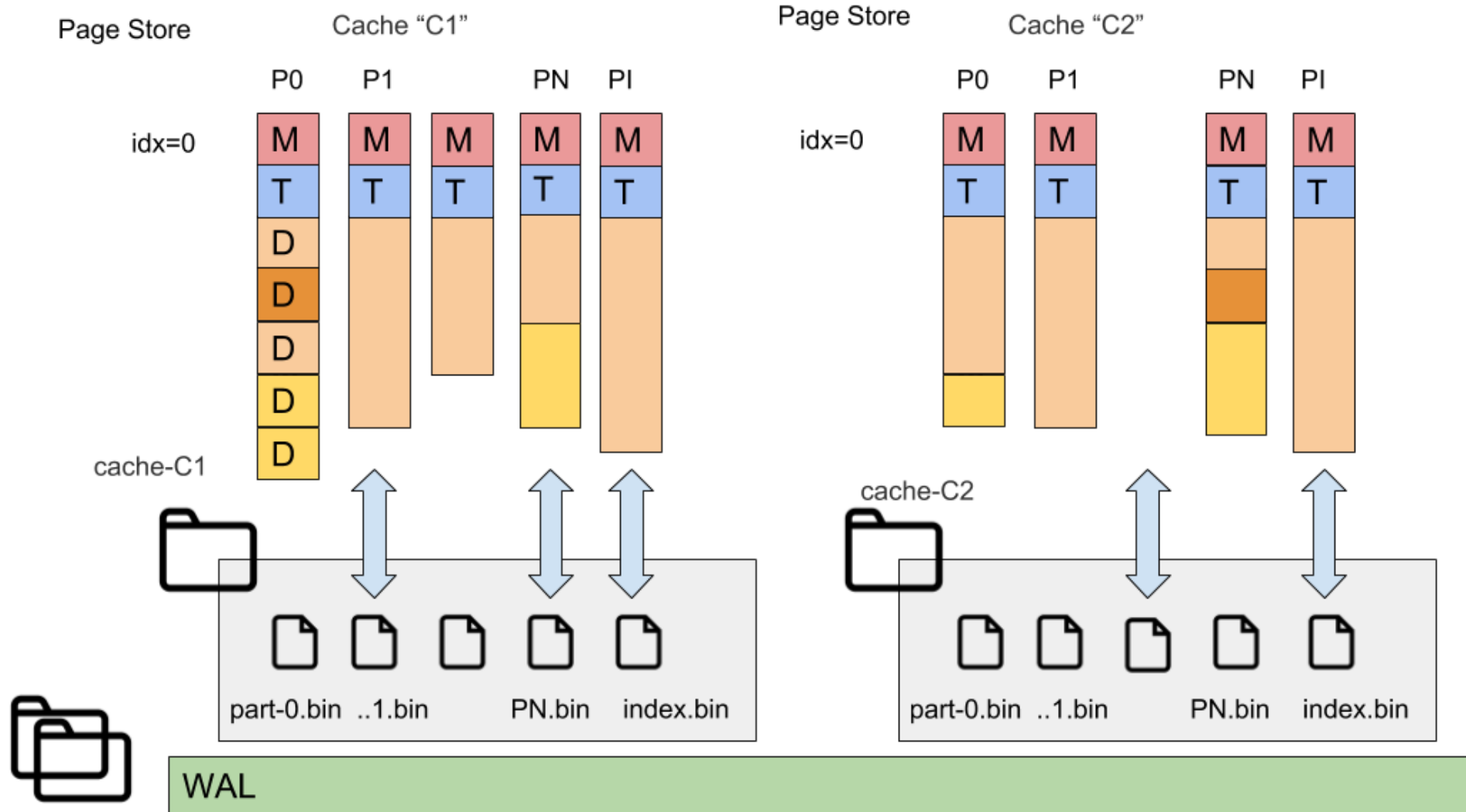Page store
(main storage for caches Data)

# File Structure

- File per Partition
- Folder per cache
- WAL shared by all caches
- Indexes are shared by all partitions

# Page store: File Structure

# Locate page placement: File Offset=Index*pageSize
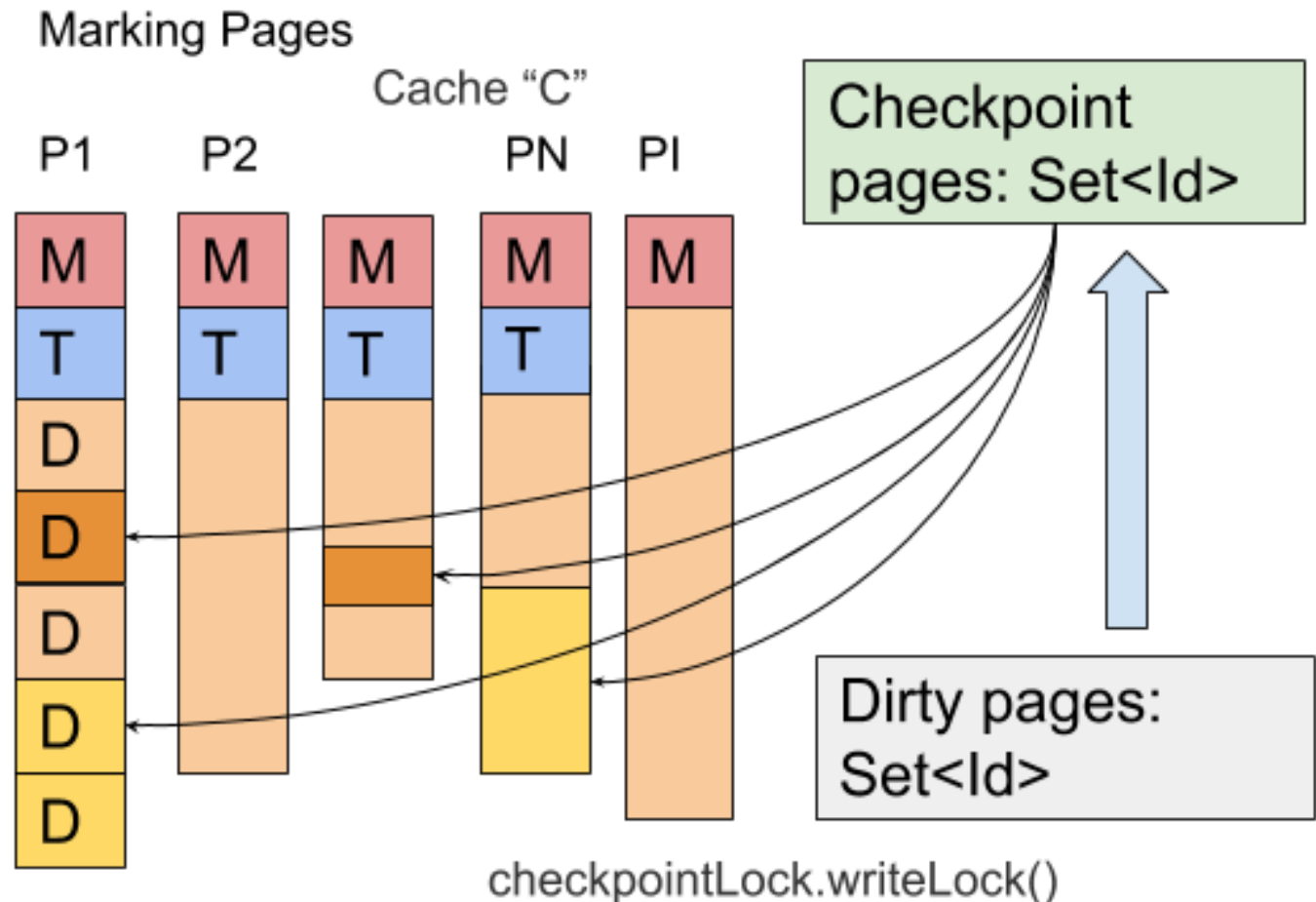
# Checkpointing

# RAM to page store: checkpointing

Periodic action: timer or dirty pages percent

Fast STW collection of current dirty pages sets

It is our scope of data to be written to disk
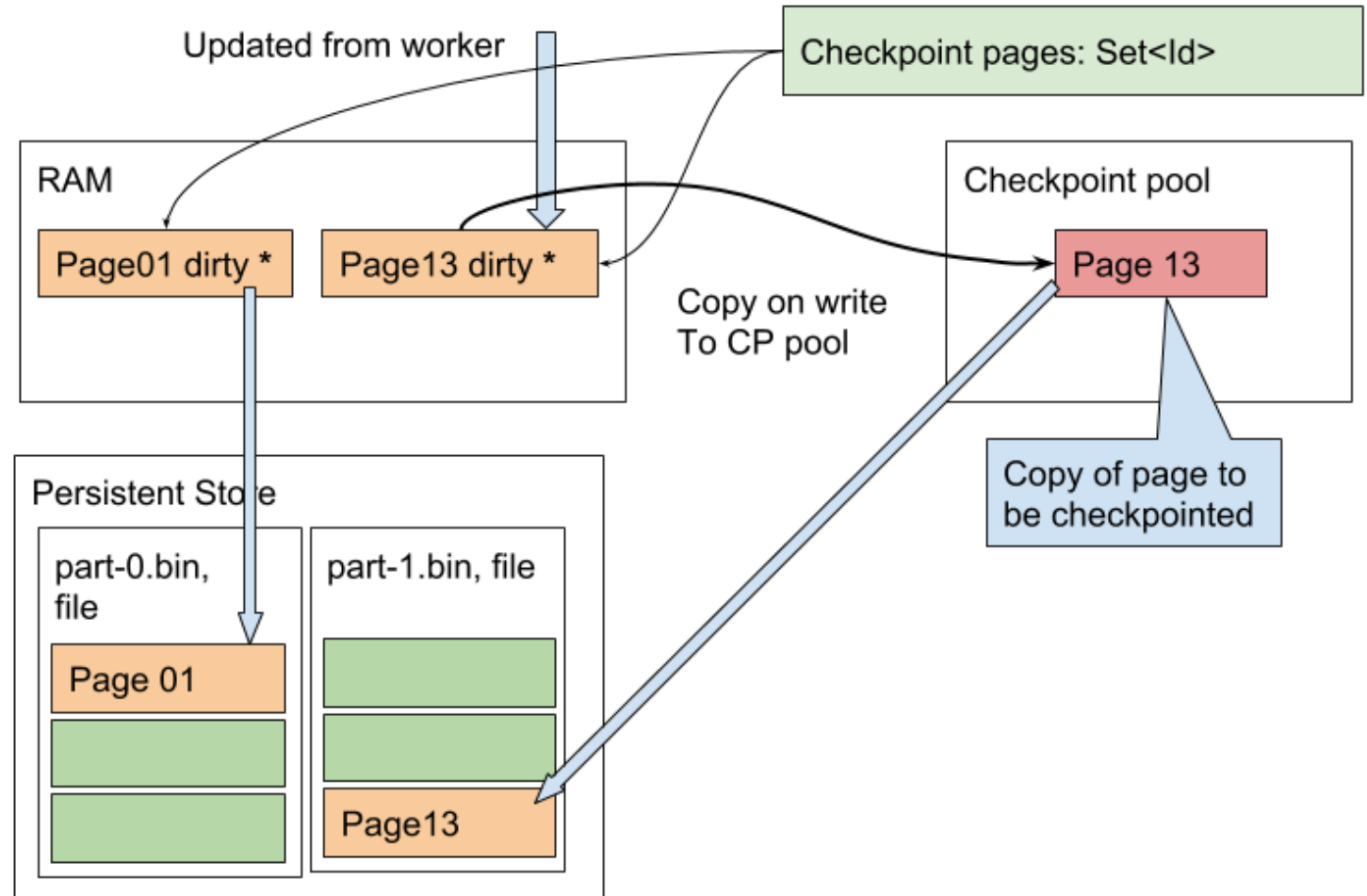
Saved page: dirty = 0



Marking Pages

Cache "C"

P1  P2  PN  PI

Checkpoint pages: Set<Id>

Dirty pages: Set<Id>

checkpointLock.writeLock()

# Process of writing

Page conflict during checkpoint => Copy on write.

CP buffer/pool

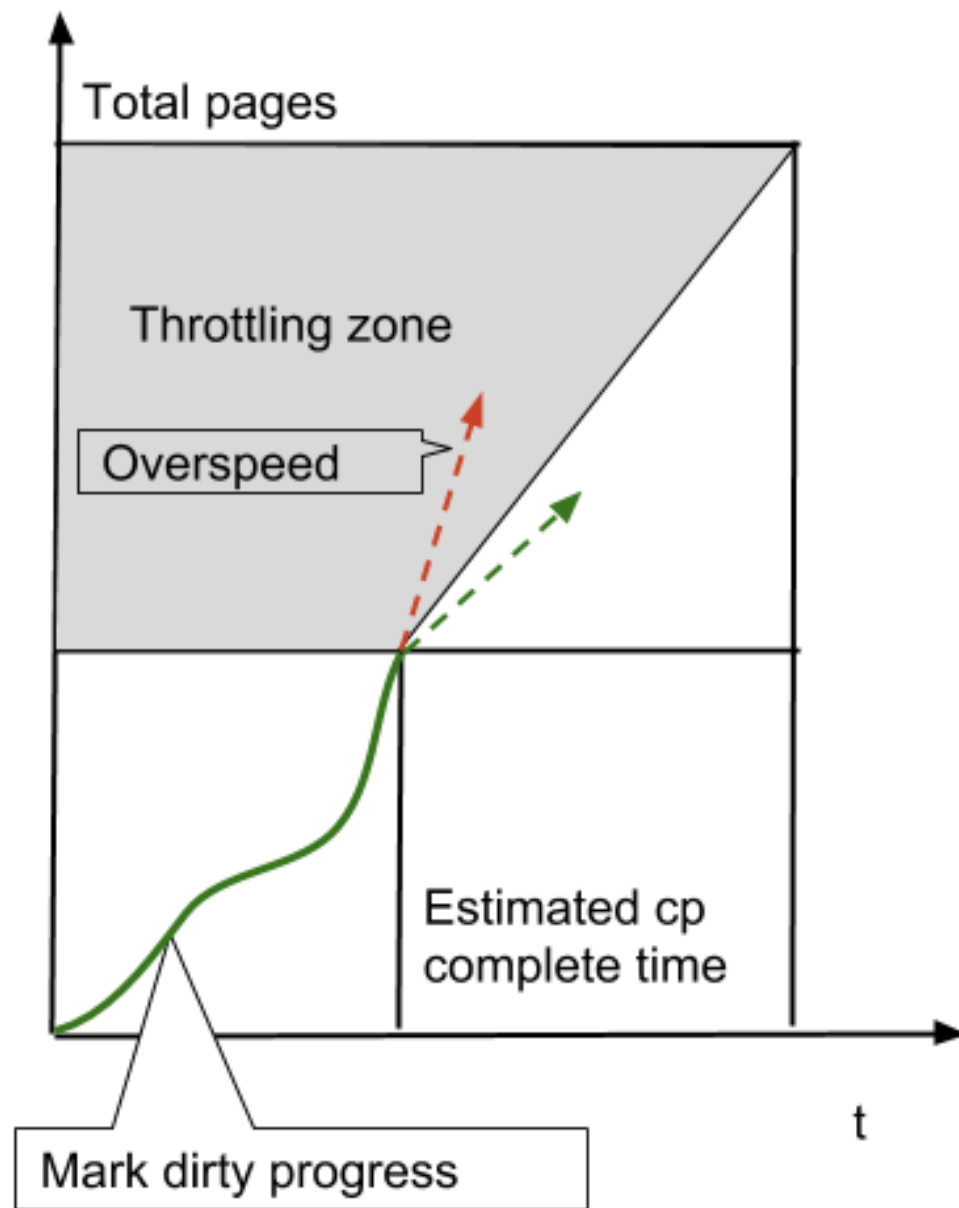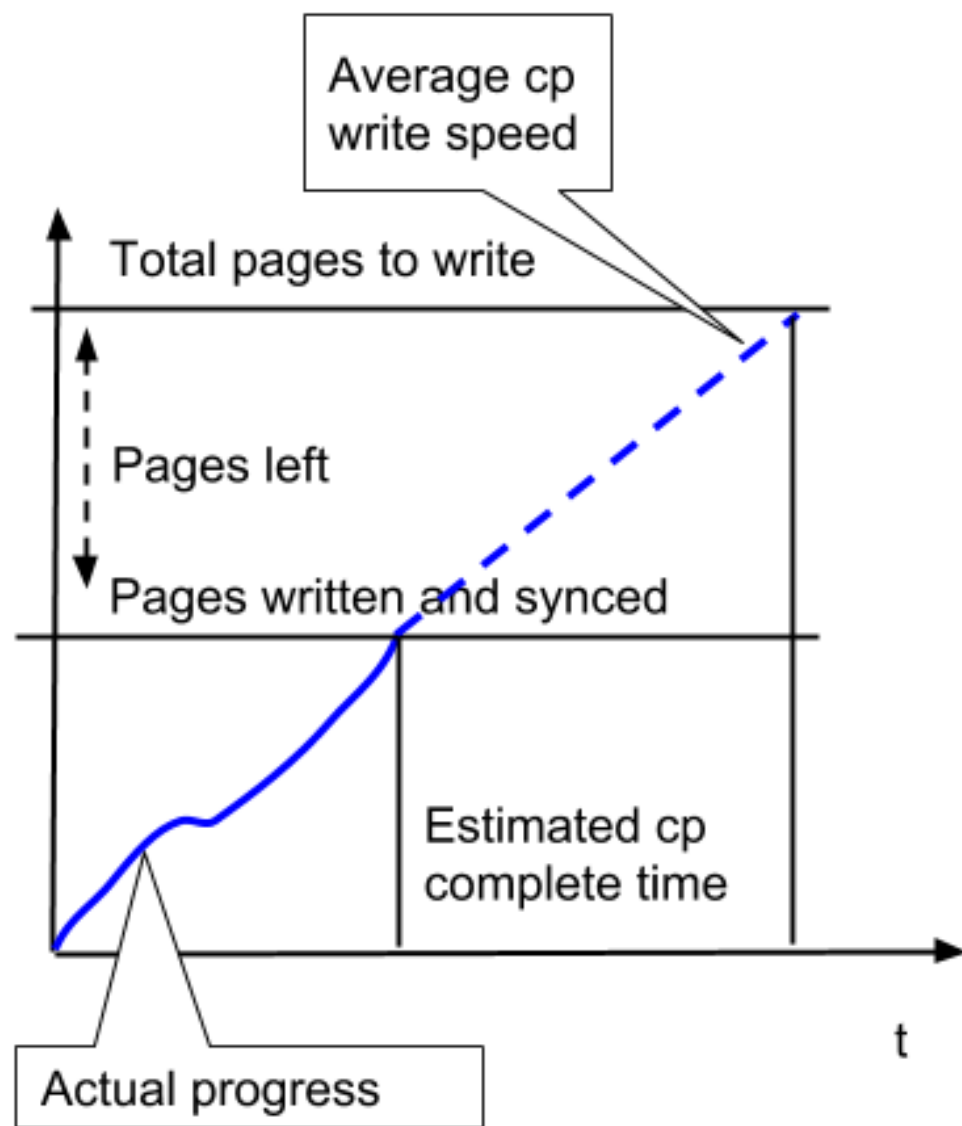Overflow protected by exponential back off always enabled

# AGENDA

What is Apache Ignite and caches
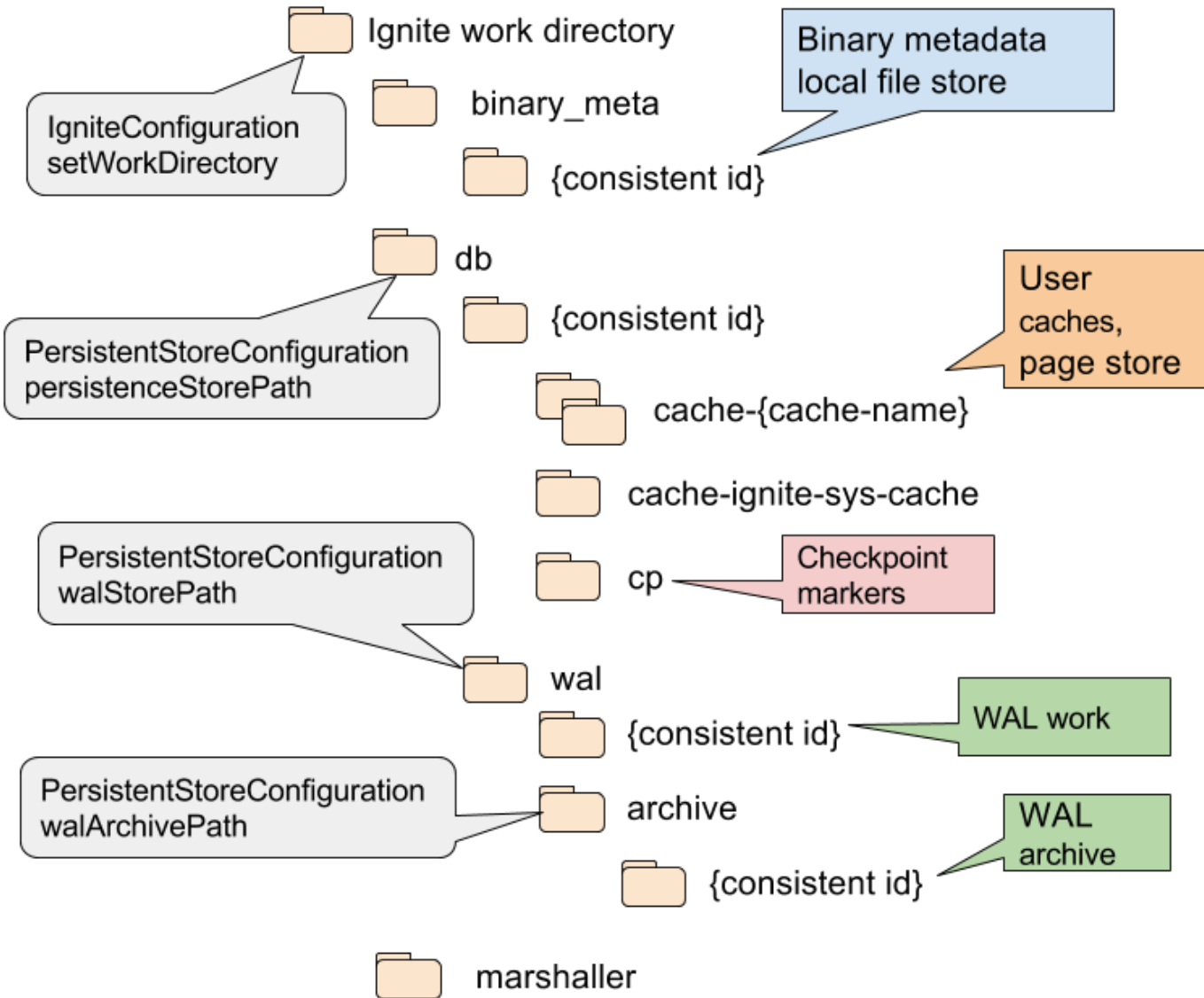
How to find an ideal match node

How to find some entry and update it

How the storage and redo log are organized

➕ Over speed protection ➕

Average cp write speed

Total pages to write

Pages left

Pages written and synced

Estimated cp complete time

Actual progress

t

Total pages

Throttling zone

Overspeed

Estimated cp complete time

Mark dirty progress

t

# Files layout

• Consistent Id – randomly or user specified node ID

Not covered

• Marshaller cache

• Binary Metadata

# Summary

# Do and Don't of today

- Don't write odd staff to your DB

- Length in bytes still matters

- Separate WAL & Page Store

- Don't set several nodes to share one HDD

- Use SSD where possible

**+** <u>**Main don't of today**</u> **+**

**Don't write your own database**

**(if you still want hardcore, join Apache Ignite community**
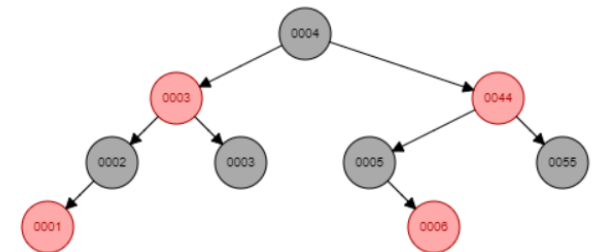**dev@ignite.apache.org)**

<u>**https://ignite.apache.org/community/contribute.html**</u>

# Links

http://ignite.apache.org/

https://apacheignite.readme.io/docs/durable-memory

https://cwiki.apache.org/confluence/display/IGNITE/Ignite+Durable+Memory+-+under+the+hood

https://apacheignite.readme.io/docs/distributed-persistent-store

https://cwiki.apache.org/confluence/display/IGNITE/Ignite+Persistent+Store+-+under+the+hood

https://www.cs.usfca.edu/~galles/visualization/Algorithms.html

# Should you have any question

dpavlov@apache.org

or just google

"Dmitriy Pavlov Apache"

Apache Ignite

How to use

user@ignite.apache.org


Related to contribution

dev@ignite.apache.org

Ignite
<summit>

May, 25, Online and free

HL HighLoad++
Весна 2021